# Sensitivity-based Investigation of Threshold Voltage Variability in 32-nm Flash Memory Cells

*Valentina Bonfiglio[1], Giuseppe Iannaccone[1,2]*

[1]Dipartimento di Ingegneria dell'Informazione and [2]SEED Center, PUSL, Università di Pisa.
Email: {valentina.bonfiglio, g.iannaccone}@iet.unipi.it

*Abstract—* **We investigate variability of a 32 nm flash memory cell with a methodology based on sensitivity analysis performed with a limited number of TCAD simulations. We show that - as far as the standard deviation of the threshold voltage is concerned - our method provides results in very good agreement with those from three-dimensional atomistic statistical simulations, with a computational burden that is orders of magnitude smaller. We show that the proposed approach is a powerful tool to understand the role of the main variability sources and to explore the device design parameter space.**

## INTRODUCTION

Non-volatile memory fabrication processes undergo even more aggressive scaling than CMOS technology for logic applications, as a means to increase bit density in response to the evolving demands of multimedia applications and mass storage. This exacerbates the device variability issue, which is especially acute in the case of multi-bit cells, where only few tens of electrons in the floating gate can separate two different logic levels [1].

The problem is particularly severe because floating gate cells must be designed and characterized for more than eight standard deviations, and therefore the second order moment of the probability distribution is hardly sufficient. [2]

In this paper we show that a recently proposed TCAD-based sensitivity analysis [3], can provide very interesting results at a small computational cost, at least for the calculation of the standard deviation of the threshold voltage. In the framework of the ENIAC Joint Undertaking MODERN project [4], we have considered a template device structure for a 32 nm CMOS flash memory cell, for which variability assessments based on three-dimensional atomistic statistical simulations and the impedance field method have been published [5]. We analyze the impact of variability sources such as random dopant distribution (RDD) [6], line-edge roughness (LER), line-width roughness (LWR), [7-8] interface trapped charge (ITC) [9], oxide thickness fluctuations (OTF) [10].

The template device structure is illustrated in Figure 1. It is a simplified polisilicon floating gate device with dimensions typical of a 32 nm technology (indicated in the table in Figure 1), generated at the crossing point of two orthogonal lines of width 32 nm. Control gate and floating gate consist of polysilicon and are separated by an ONO (oxide-nitride-onide) layer of 4-3-5 nm. The tunnel oxide thickness is 8 nm. Substrate is boron doped ($2 \times 10^{18}\,cm^{-3}$), and arsenic doping of source and drain is symmetric with a maximum of $10^{20}\,cm^{-3}$, Gaussian shape, and junction depth of 25 nm. Additional details are available in Ref. [5].



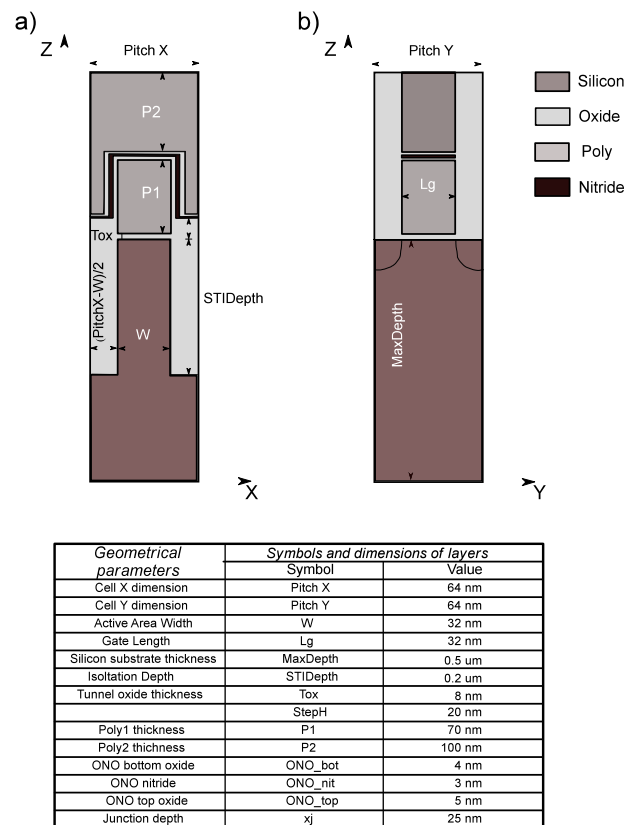| Geometrical parameters | Symbols and dimensions of layers | |
|---|---|---|
| | Symbol | Value |
| Cell X dimension | Pitch X | 64 nm |
| Cell Y dimension | Pitch Y | 64 nm |
| Active Area Width | W | 32 nm |
| Gate Length | Lg | 32 nm |
| Silicon substrate thickness | MaxDepth | 0.5 um |
| Isolation Depth | STIDepth | 0.2 um |
| Tunnel oxide thickness | Tox | 8 nm |
| | StepH | 20 nm |
| Poly1 thickness | P1 | 70 nm |
| Poly2 thichness | P2 | 100 nm |
| ONO bottom oxide | ONO_bot | 4 nm |
| ONO nitride | ONO_nit | 3 nm |
| ONO top oxide | ONO_top | 5 nm |
| Junction depth | xj | 25 nm |

Figure 1: Device structure and geometrical parameters of the template 32 nm flash memory under investigation.

## METHODOLOGY

The approach proposed is described in detail in [3]. First, all process and geometry variability causes are expressed in terms of a set of synthetic independent variability sources. Then, TCAD-based sensitivity analysis is used to evaluate the contribution to the dispersion of electrical parameters (e.g. the threshold voltage $V_{th}$) of each independent source. This step is based on the assumption that the effect of each source is sufficiently small that first-order linearization is applicable. Also in the case of the 32 nm Flash memory [5], the variance of the threshold voltage due to combined effect computed with 3D atomistic statistical is shown to be very close to the sum of the variances due to individual effects, giving us confidence in the linear approximation.

As an example, let us consider the case of LER, considering the illustration in Fig. 2, where the 32 nm device is shown with the $y$ axis running along the channel length direction, the $x$ axis perpendicular to the device plane and the $z$ axis running along the channel width.

We can translate line edge roughness in terms of the dispersion of the average position of both gate edges along the y axis ( $y_1$ and $y_2$, where $\langle y_1 \rangle = 0$ and $\langle y_2 \rangle = L$). This in turn translates into gate length dispersion. We assume that parameters $y_1, y_2$ are only affected by LER and are physically independent. The average edge position is a random function $g(z)$ with zero mean value and Gaussian autocorrelation $r(d) \equiv \langle g(z)g(z+d) \rangle$ characterized by correlation length $\Lambda_L$ and mean square amplitude $\Delta_L$, i.e.:

$$r(d) = \Delta_L^2 e^{\frac{d^2}{2\Lambda_L^2}} \qquad (1)$$

from which we can write the variance of $g$ as

$$\sigma_g^2 \equiv \langle g^2 \rangle = \frac{1}{W^2} \left\langle \int_0^W g(z_1)dz_1 \cdot \int_0^W g(z_2)dz_2 \right\rangle. \qquad (2)$$
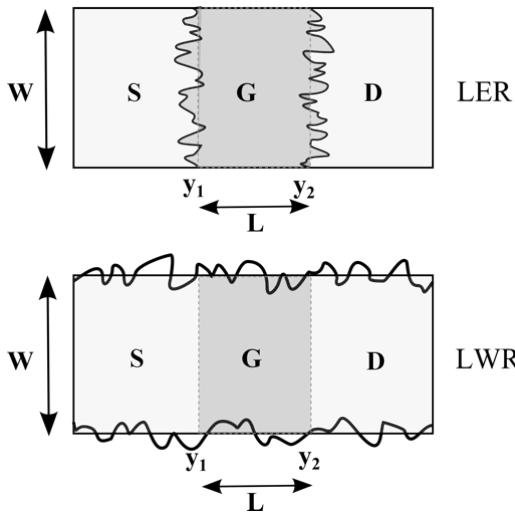


Figure 2: Illustration of the approach to the evaluation of line edge roughness (above) and line-width roughness (below).

If we compute (2) considering (1) we find:

$$\sigma_{LER}^2 = \frac{2\Delta_L^2 \Lambda_L}{W^2}\left[ \Lambda_L \left( e^{-\frac{W^2}{2\Lambda_L^2}} - 1 \right) + \sqrt{\frac{\pi}{2}} W \mathrm{erf}\left( \frac{W}{\sqrt{2}\Lambda_L} \right) \right] \qquad (3)$$

The variance of $V_{th}$ due to line edge roughness is:

$$\sigma_{V_{th}LER}^2 = \left( \frac{\partial V_{th}}{\partial y_1} \right)^2 \sigma_{y_1}^2 + \left( \frac{\partial V_{th}}{\partial y_2} \right)^2 \sigma_{y_2}^2 = 2\left( \frac{\partial V_{th}}{\partial L} \right)^2 \sigma_{LER}^2, (4)$$

where $y_1, y_2$ in (4) are the average gate edges indicated in Fig. 2. All required derivatives can be computed with TCAD sensitivity analysis as illustrated in Fig. 3 (left). The very same approach can be used for LWR.

In the case of OTF we must consider surface roughness with a two dimensional Gaussian autocorrelation

$$r(x_a,y_a,x_b,y_b) = \Delta_S^2 \exp\left( -\frac{(x_b - x_a)^2 + (y_b - y_a)^2}{2\Lambda_S^2} \right), \qquad (5)$$

characterized by correlation length $\Lambda_S$ and mean square amplitude $\Delta_S$, which corresponds to a variance of the average position of the interface:

$$\sigma_{SR}^2 = \frac{2\pi\Lambda_S^2\Delta_S^2}{L^2 W^2}\left[ L \cdot \mathrm{erf}\left( \frac{L}{\sqrt{2}\Lambda_S} \right) + \sqrt{\frac{2}{\pi}}\Lambda_S \left( e^{-\frac{L^2}{2\Lambda_S^2}} - 1 \right) \right]$$
$$\times \left[ W \mathrm{erf}\left( \frac{W}{\sqrt{2}\Lambda_S} \right) + \sqrt{\frac{2}{\pi}}\Lambda_S \left( e^{-\frac{W^2}{2\Lambda_S^2}} - 1 \right) \right], \qquad (6)$$

The variance of the threshold voltage due to OTF is therefore

$$\sigma_{V_{th}SR}^2 = \sum_m \left( \frac{\partial V_{th}}{\partial s_m} \right)^2 \sigma_{OTF}^2, \qquad (7)$$

where $s_m$ are all positions of the interfaces between dielectric layers and between dielectric and conducting or semiconducting layers. Also in this case, all derivatives can be computed with TCAD simulations following the example of Figure 3 (right).

For LER and LWR, we consider a Gaussian autocorrelation with mean square amplitude $\Delta_L = 1.5$ nm and correlation length $\Lambda_L = 20$ nm. For OTF, we consider a Gaussian autocorrelation with with mean square amplitude $\Delta_S = 0.2$ nm and correlation length $\Lambda_S = 18$ nm.

Results are compared in Table I with those obtained from 3D atomistic simulations on 1000 samples performed with GARAND [5], in which the same statistical properties have been considered for LER, LWR, and OTF. The obtained standard deviation are practically identical. As in [5], the threshold voltage is defined with a current criterion of 100 nA for a drain-to-source voltage of 100 mV.
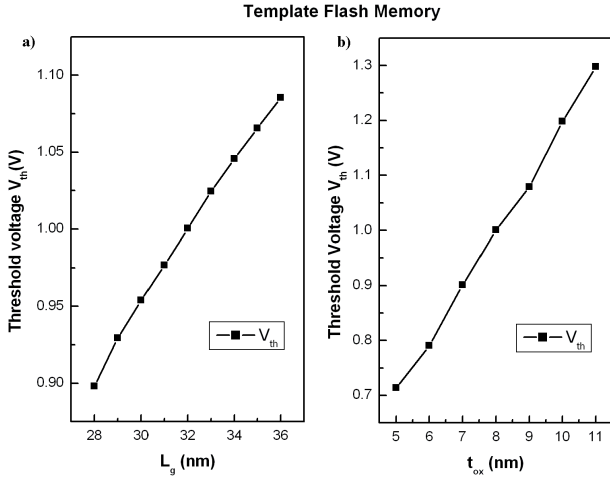
**Template Flash Memory**

a) threshold voltage $V_{th}$ (V) vs $L_g$ (nm)

b) Threshold Voltage $V_{th}$ (V) vs $t_{ox}$ (nm)

Figure 3 a) Threshold voltage as a function of gatelength Lg and b) threshold voltage as a function of tunnel oxide thickness $t_{ox}$ for the template Flash Memory as computed from TCAD simulations.

For random discrete dopants (RDD) [6] and interface trapped charge (ITC) [10], we adopt an approach based on a propagator with a very coarse granularity, which is in principle very close to the concept of impedance field method [11]. As a difference with respect to the situation already described in [4], we here have to perform 3D simulations, since the Flash memory cells cannot be reduced to 2D structures.

For a given variation of doping concentration $\Delta N_A(x,y,z)$ with respect to the nominal value we can write the following expression for the variation of $V_{th}$:

$$\Delta V_{th} = \int K(x,y,z)\Delta N_A(x,y,z)dxdydz \qquad (8)$$

where $K(x,y,z)$ has the role of a propagator. The expression requires the linearity assumption to hold.

To conveniently compute the propagator $K$, we can assume that $K$ is a smooth function of $x$, $y$, and $z$, and move from the continuum to a discrete space, partitioning the active area in small boxes. Now we can write:

$$\Delta V_{th} = \sum_i \Delta V_{th_i} = \sum_i K_i \Delta N_i \qquad (9)$$

The sum runs over all boxes, $\Delta N_i$ is the variation of the number of dopants in box $i$, and $\Delta V_{th_i}$ is the threshold voltage variation if only dopants in box $i$ are varied.

In practice, we multiply doping in box $i$ by a factor $(1+\alpha)$ and compute $\Delta V_{thi}$ with TCAD simulations. Therefore we have

$$\Delta N_i = \alpha N_i$$
$$\Delta V_{th_i} = \alpha K_i N_i \qquad (10)$$

so that (9) becomes,

$$\Delta V_{th} = \sum_i \left(\frac{\Delta V_{th_i}}{\alpha}\right)\alpha = \sum_i \left(\frac{\Delta V_{th_i}}{\alpha}\right)\frac{\Delta N_i}{N_i} \qquad (11)$$

If we finally assume that doping variations in different boxes are independent Poisson processes, we can write

$$\sigma_{V_{thRDD}}^2 = \sum_i \left(\frac{\Delta V_{th_i}}{\alpha}\right)^2 \frac{1}{N_i} = \sum_i \sigma_{V_{thRDD}}^{2\,[i]}, \qquad (12)$$

The threshold voltage dispersion due to RDD only requires a single TCAD simulation for each box, and an integral of the doping profile in each box. To evaluate the most convenient level of granularity in device partitioning, we have made tests with different box sizes, as reported in the table in Figure 4.



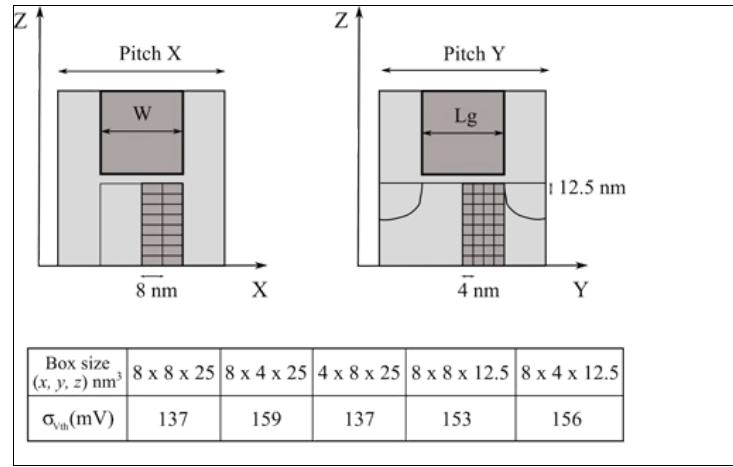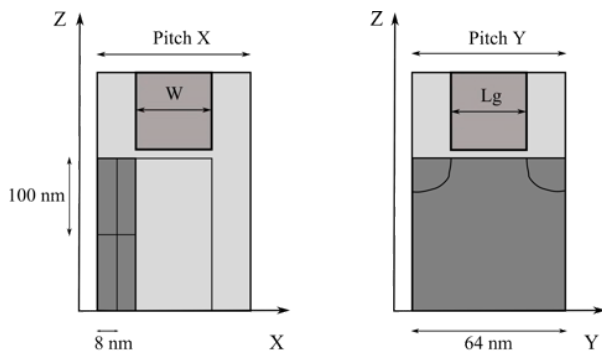| Box size (x, y, z) nm³ | 8 x 8 x 25 | 8 x 4 x 25 | 4 x 8 x 25 | 8 x 8 x 12.5 | 8 x 4 x 12.5 |
|---|---|---|---|---|---|
| $\sigma_{Vth}$(mV) | 137 | 159 | 137 | 153 | 156 |

Figure 4 above: transversal (left) and longitudinal (right) device cross sections for the assessment of the proper box partitioning. Below: computed standard deviation of the threshold voltage as a function of the box size for different choices of the partition.

We have evaluated that a partition of the three dimensional silicon body in 64 boxes of size $8\times4\times12.5$ nm³ represents a good trade-off between computing time and accuracy. Considering that we can exploit the symmetry of the structure also along the transport direction at very low drain-to-source voltage, only sensitivities corresponding to 32 boxes must be computed with TCAD simulations.

For ITC, the situation is similar: we assume an average trap density of $5\times10^{11}$ cm⁻² and partition the tunnel oxide in tales of $100\times8\times64$ nm³, for a total of only four simulations, if the symmetry of the nominal structure is exploited. As can be seen in Figure 5, finer partitions do not lead to a different estimation of the threshold voltage dispersion.

| Box size $(x, y, z)$ nm$^3$ | 16 x 64 x 200 | 8 x 64 x 100 | 16 x 64 x 50 | 8 x 32 x 50 | 4 x 32 x 50 |
|---|---|---|---|---|---|
| $\sigma_{V_{th}}$(mV) | 27 | 59 | 56 | 59 | 59 |

Figure 5: Region partitioning in boxes $100 \times 8 \times 64$ nm$^3$ for the evaluation of propagators due to interface trapped charge. Left: transversal cross section. Right: longitudinal cross section.

The effect of RDD and ITC on the threshold voltage have been compared in Table 1 with direct simulation of a statistical ensemble done at the University of Glasgow through GARAND [5] obtained simulating samples of 1000 microscopically different devices. Considering that statistical simulations have been performed on ensembles of $N$=1000 devices, the mean square relative error on the estimated standard deviation of the threshold voltage is $(2N)^{-0.5}$, i.e., 2.2%: all terms lie within or very close to the error bars of statistical simulations.

TABLE 1 STANDARD DEVIATION OF THE THRESHOLD VOLTAGE DUE TO LER AND LWR OBTAINED WITH THE METHOD PROPOSED IN [3] AND WITH STATISTICAL SIMULATION IN [5].

| $\sigma_{V_{th}}$ (mV) | Our method [3] | Atomistic Sim. [5] |
|---|---|---|
| LER | 46 | 48 |
| LWR | 28 | 26 |
| OTF | 14 | 14 |
| RDD | 156 | 144 |
| ITC | 59 | 67 |

## CONCLUSION

We have proposed a methodology for the quantitative evaluation of the effects of the main mechanisms affecting threshold voltage variability, based on the careful identification of the main independent and relevant physical quantities. Our approach requires the calculation of partial derivatives of $V_{th}$ with respect to device structure parameters, that can be obtained with a very limited number of TCAD simulations. We have shown that in all cases we are able to obtain results in good agreement with 3D atomistic statistical simulations [5] at a much smaller computational cost. We qualify this statement to the second order moment of the threshold voltage distribution,

because the proposed approach does not provide information on the far tails of the distribution, which are important for large Flash memory arrays, and would require extension of the method to higher order terms.

Our approach has some advantages over statistical modeling, not only because is orders of magnitude faster, but also because it represents a powerful tool for understanding the impact of individual factors and to efficiently explore the design space using tools already available and routinely used by technology developers .

### REFERENCES

[1] A. Calderoni, P. Fantini, A. Ghetti, A. Marmiroli, "Vth fluctuations in nanoscale floating gate memories, Proc. SISPAD, Sept. 9-11, 2008, pp. 49-52.

[2] A. Spessot, A. Calderoni, P. Fantini, A. S. Spinelli, C. Monzio Compagnoni, F. Farina, A. L. Lacaita, and A. Marmiroli, "Variability effects on the VT distribution of nanoscale NAND Flash memories," in Proc. IRPS, 2010, pp. 970–974.

[3] V. Bonfiglio and G. Iannaccone, "An Approach Based on Sensitivity Analysis for the Evaluation of Process Variability in Nanoscale MOSFETs," IEEE Transactions on Electron Devices, vol. 58, no. 8, pp. 2266-2273, 2011.

[4] Deliverable D.2.2.3 of the ENIAC MODERN Project, 2011.(Project website: www.eniac-modern.org).

[5] G. Roy, A. Ghetti, A. Benvenuti, A. Erlebach, and A. Asenov, "Comparative Simulation Study of the Different Sources of Statistical Variability in Contemporary Floating-Gate Nonvolatile Memory," IEEE Transactions on Electron Devices, vol. 58, no. 12, pp. 4155-4163, Dec. 2011.

[6] H.-S. Wong, Y. Taur, "Three-dimensional "atomistic" simulation of discrete random dopant distribution effects in sub-0.1 mm MOSFETs", Tech. Dig. IEDM 1993, pp. 705-708, 1993.

[7] A. Asenov, S. Kaya, and A. R. Brown, "Intrinsic parameter fluctuations in decananometre MOSFETs introduced by gate line edge roughness," IEEE Trans. Electron Devices, vol. 50, no. 5, pp. 1254–1260, May 2003.

[8] Ji-Young Lee, Jangho Shin, Hyun-Woo Kim, Sang-Gyun Woo, Han-Ku Cho, Woo-Sung Han, and Joo-Tae Moon, "Effect of line-edge roughness (LER) and line-width roughness (LWR) on sub-100-nm device performance", Proc. SPIE 5376, 426 (2004).

[9] A. Asenov, S. Kaya and J. H. Davies, "Intrinsic threshold voltage fluctuations in decanano MOSFETs due to local oxide thickness variations," IEEE Transactions on Electron Devices, vol. 49, pp. 112–119, 2002.

[10] C. L. Alexander, A. R. Brown, J. R. Watling, and A. Asenov, "Impact of single charge trapping in nano-MOSFETs—Electrostatics versus transport effects," IEEE Trans. Nanotechnol., vol. 4, no. 3, pp. 339–344,May 2005.

[11] W. Shockley, J. A. Copeland, R. P. James, "The impedance field method of noise calculation in active semiconductor devices", in Quantum Theory of Atoms, Molecules, and the Solid State, A tribute to John C. Slater. Edited by Per-Olov Loewdin. New York: Academic Press, 1966, p.537.